

GAUSSIAN CONDITIONAL RANDOM FIELDS FOR CLASSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, a Gaussian conditional random field model for structured binary classification (GCRFBC) is proposed. The model is applicable to classification problems with undirected graphs, intractable for standard classification CRFs. The model representation of GCRFBC is extended by latent variables which yield some appealing properties. Thanks to the GCRF latent structure, the model becomes tractable, efficient, and open to improvements previously applied to GCRF regression. Two different forms of the algorithm are presented: GCRF-BCb (GCRGBC - Bayesian) and GCRFBCnb (GCRFBC - non-Bayesian). The extended method of local variational approximation of sigmoid function is used for solving empirical Bayes in GCRFBCb variant, whereas MAP value of latent variables is the basis for learning and inference in the GCRFBCnb variant. The inference in GCRFBCb is solved by Newton-Cotes formulas for one-dimensional integration. Both models are evaluated on synthetic data and real-world data. It was shown that both models achieve better prediction performance than relevant baselines. Advantages and disadvantages of the proposed models are discussed.

1 INTRODUCTION

Increased quantity and variety of sources of data with correlated outputs, so called structured data, created an opportunity for exploiting additional information between dependent outputs to achieve better prediction performance. One of the most successful probabilistic models for structured output classification problems are conditional random fields (CRF) (Sutton & McCallum, 2006). The main advantages of CRFs lie in their discriminatory nature, resulting in the relaxation of independence assumptions and the label bias problem that are present in many graphical models. Aside of many advantages, CRFs also have many drawbacks mostly resulting in high computational cost or intractability of inference and learning. A wide range of different approaches of tackling these problems has been proposed, and they motivate our work, too.

One of the popular methods for structured regression based on CRFs – Gaussian conditional random fields (GCRF) – has the form of multivariate Gaussian distribution (Radosavljevic et al., 2010). The main assumption of the model is that the relations between outputs are presented in quadratic form. It has convex loss function and, consequently, efficient inference and learning, and expensive sampling methods are not used.

In this paper, a new model of Gaussian conditional random fields for binary classification is proposed (GCRFBC). GCRFBC builds upon regression GCRF model which is used to define latent variables over which output dependencies are defined. The model assumes that discrete outputs y_i are conditionally independent conditioned on continuous latent variables z_i which follow a distribution modeled by a GCRF. That way, relations between discrete outputs are not expressed directly. Two different inference and learning approaches are proposed in this paper. The first one is based on evaluating empirical Bayes by marginalizing latent variables (GCRFBCb), whereas MAP value of latent variables is the basis for learning and inference in the second model (GCRFBCnb). In order to derive GCRFBCb model and its learning procedure the variational approximation of Bayesian logistic regression (Jaakkola & Jordan, 2000) is generalized.

Compared to CRFs and structured SVM classifiers, the GCRFBC models have some appealing properties:

- The model is applicable to classification problems with undirected graphs, intractable for standard classification CRFs. Thanks to the GCRF latent structure, the model becomes tractable, efficient and open to improvements previously applied to GCRF regression models.
- Defining correlations directly between discrete outputs may introduce unnecessary noise to the model (Tan et al., 2010). This problem can be solved by defining structured relations on a latent continuous variable space.
- In case that unstructured predictors are unreliable, which is signaled by their large variance (diagonal elements in the covariance matrix), it is simple to marginalize over latent variable space and obtain better results.

GCRFBC model is relying on the assumption that the underlying distribution of latent variables is multivariate normal distribution, due to that in the case when this distribution cannot be fitted well to the data (e.g. when the distribution of latent variables is multimodal) the model will not perform as well as it is expected. The proposed models are experimentally tested on both synthetic and real-world datasets in terms of predictive performance and computation time. In experiments with synthetic datasets, the results clearly indicate that the the empirical Bayes approach (GCRFBCb) better exploits output dependence structure, more so as the variance of the latent variables increases. We also tested both approaches on real-world datasets of predicting ski lift congestion, gene function classification, classification of music according to emotion and highway congestion. Both GCRFBC models outperformed ridge logistic regression, lasso logistic regression, neural network, random forest, and structured SVM classifiers, demonstrating that the proposed models can exploit output dependencies in a real-world setting.

2 RELATED WORK

An extensive review of binary and multi-label classification with structured output is provided in Su (2015). A number of different studies related to graph based methods for regression can be found in the literature (Fox, 2015). CRFs were successfully applied on a variety of different structured tasks (Cotterell & Duh, 2017; Zhang et al., 2015; Masada & Bunescu, 2017; Zia et al., 2018) and different model adaptations can be found in literature Kim (2017); Maaten et al. (2011). Recently, successful unifications of deep learning and CRFs have been proposed Chen et al. (2016); Kosov et al. (2018). Moreover, implementation of deep neural networks as potential functions is presented in form of structure prediction energy networks (SPEN) Belanger & McCallum (2016); Belanger et al. (2017). Adaptation of normalizing flows in SPEN structure is presented in Lu & Huang (2019).

An extensive review on topic of binary and multi-label classification with structured output is provided in Su (2015). Large number of different studies related to graph based methods for regression can be found in the literature (Fox, 2015). CRFs were successfully applied on a variety of different structured tasks, such as: low-resource named entity recognition (Cotterell & Duh, 2017), image segmentation (Zhang et al., 2015), chord recognition (Masada & Bunescu, 2017) and word segmentation (Zia et al., 2018). The mixture of CRFs capable to model data that come from multiple different sources or domains is presented in Kim (2017). The method is related to the well known hidden-unit CRF (HUCRF) (Maaten et al., 2011). The conditional likelihood and expectation minimization (EM) procedure for learning have been derived there. The mixtures of CRF models were implemented on several real-world applications resulting in prediction improvement. Recently, a model based on unification of deep learning and CRF was developed by Chen et al. (2016). The deep CRF model showed better performance compared to either shallow CRFs or deep learning methods on their own. Similarly, the combination of CRFs and deep convolutional neural networks was evaluated on an example of environmental microorganisms labeling (Kosov et al., 2018). The spatial relations among outputs were taken in consideration and experimental results have shown satisfactory results.

The GCRF model was first implemented for the task of low-level computer vision (Tappen et al., 2007). Since then, various different adaptations and approximations of GCRF were proposed (Radostavljevic et al., 2014). The parameter space for the GCRF model is extended to facilitate joint modelling of positive and negative influences (Glass et al., 2016). In addition, the model is extended by bias term into link weight and solved as a part of convex optimization. Semi-supervised

marginalized Gaussian conditional random fields (MGCRF) model for dealing with missing variables was proposed by Stojanovic et al. (2015). The benefits of the model were proved on partially observed data and showed better prediction performance than alternative semi-supervised structured models. A comprehensive review of continuous conditional random fields (CCRF) was provided in Radosavljevic et al. (2010). The sparse conditional random fields obtained by l_1 regularization are first proposed and evaluated by Wytock & Kolter (2013). Additionally, Frot et al. (2018) presented GCRF with the latent variable decomposition and derived convergence bounds for the estimator that is well behaved in high dimensional regime. An adaptation of GCRF on discrete output was briefly discussed in Radosavljevic (2011), as a part of future work. This discussion motivates our work, but our approach is different in technical aspects.

3 METHODOLOGY

In this section we first present already known GCRF model for regression and then we propose GCRFBC model for binary classification and two approaches to inference and learning.

3.1 BACKGROUND MATERIAL

GCRF is a discriminative graph-based regression model (Radosavljevic et al., 2010). Nodes of the graph are variables $\mathbf{y} = (y_1, y_2, \dots, y_N)$, which need to be predicted given a set of features \mathbf{x} . The attributes $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$ interact with each node y_i independently of one another, while the relations between outputs are expressed by pairwise interaction function. In order to learn parameters of the model, a training set of vectors of attributes x and real-valued response variables y are provided. The generalized form of the conditional distribution $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ is:

$$P(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{Z(\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})} \exp \left(- \sum_{i=1}^N \sum_{k=1}^K \alpha_k (y_i - R_k(\mathbf{x}_i))^2 - \sum_{i \neq j} \sum_{l=1}^L \beta_l S_{ij}^l (y_i - y_j)^2 \right) \quad (1)$$

First sum models relations between outputs y_i and corresponding input vector \mathbf{x}_i and the second one models pairwise relations between nodes. $R_k(\mathbf{x}_i)$ represents an unstructured predictor of y_i for each node in the graph and S_{ij}^l is value that expresses similarity between nodes i and j in graph l . Unstructured predictor can be any regression model that gives prediction of output y_i for given attributes \mathbf{x}_i . K is the total number of unstructured predictors. L is the total number of graphs (similarity functions). Graphs can express any kind of binary relations between nodes e.g., spatial and temporal correlations between outputs. Z is a partition function and vectors $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ are learnable parameters. One of the main advantages of GCRF is the ability to express different relations between outputs by variety of graphs and ability to learn which graphs are significant for prediction. The quadratic form of interaction and association potential enables conditional distribution $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ to be expressed as multivariate Gaussian distribution (Radosavljevic et al., 2010):

$$P(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{(2\pi)^{\frac{N}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right) \quad (2)$$

Precision matrix $\Sigma^{-1} = 2Q$ and distribution mean $\boldsymbol{\mu} = \Sigma \mathbf{b}$ are defined as, respectively:

$$Q = \begin{cases} \sum_{k=1}^K \alpha_k + \sum_{h=1}^N \sum_{l=1}^L \beta_l S_{ih}^l, & \text{if } i = j \\ - \sum_{l=1}^L \beta_l S_{ij}^l, & \text{if } i \neq j \end{cases} \quad (3)$$

$$\mathbf{b}_i = 2 \left(\sum_{k=1}^K \alpha_k R_k(\mathbf{x}_i) \right) \quad (4)$$

Due to concavity of multivariate Gaussian distribution, the inference task $\underset{\mathbf{y}}{\operatorname{argmax}} P(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ is straightforward. The maximum posterior estimate of \mathbf{y} is the distribution expectation $\boldsymbol{\mu}$.

The objective of the learning task is to optimize parameters $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ by maximizing conditional log likelihood $\underset{\boldsymbol{\alpha}, \boldsymbol{\beta}}{\operatorname{argmax}} \sum_{\mathbf{y}} \log P(\mathbf{y}|\mathbf{x}, \boldsymbol{\alpha}, \boldsymbol{\beta})$. One way to ensure positive definiteness of the covariance matrix of GCRF is to require diagonal dominance (Strang et al., 1993). This can be ensured by imposing constraints that all elements of $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ be greater than 0 (Radosavljevic et al., 2010).

3.2 GCRFBC MODEL REPRESENTATION

One way of adapting GCRF to classification problem is by approximating discrete outputs by suitably defining continuous outputs. Namely, GCRF can provide dependence structure over continuous variables which can be passed through sigmoid function. That way the relationship between regression GCRF and classification GCRF is similar to the relationship between linear and logistic regression, but with dependent variables. Aside from allowing us to define a classification variant of GCRF, this may result in additional appealing properties: (i) The model is applicable to classification problems with undirected graphs, intractable for standard classification CRFs. Thanks to the GCRF latent structure, the model becomes tractable, efficient and open to improvements previously applied to GCRF regression models. (ii) Defining correlations directly between discrete outputs may introduce unnecessary noise to the model (Tan et al., 2010). We avoid this problem by defining structured relations on a latent continuous variable space. (iii) In case that unstructured predictors are unreliable, which is signaled by their large variance (diagonal elements in the covariance matrix), it is simple to marginalize over latent variable space and obtain better results.

It is assumed that y_i are discrete binary outputs and z_i are continuous latent variables assigned to each y_i . Each output y_i is conditionally independent of the others, given z_i .

The conditional probability distribution $P(y_i|z_i)$ is defined as Bernoulli distribution:

$$P(y_i|z_i) = \text{Ber}(y_i|\sigma(z_i)) = \sigma(z_i)^{y_i} (1 - \sigma(z_i))^{1-y_i} \quad (5)$$

where $\sigma(\cdot)$ is sigmoid function. Due to conditional independence assumption, the joint distribution of outputs y_i can be expressed as:

$$P(y_1, y_2, \dots, y_N|\mathbf{z}) = \prod_{i=1}^N \sigma(z_i)^{y_i} (1 - \sigma(z_i))^{1-y_i} \quad (6)$$

Furthermore, the conditional distribution $P(\mathbf{z}|\mathbf{x})$ is the same as in the classical GCRF model and has canonical form defined by multivariate Gaussian distribution. Hence, joint distribution of continuous latent variables \mathbf{z} and outputs \mathbf{y} given \mathbf{x} and $\boldsymbol{\theta} = (\alpha_1, \dots, \alpha_K, \beta_1, \dots, \beta_L)$ is the general form of the GCRFBC model defined as:

$$P(\mathbf{y}, \mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) = \prod_{i=1}^N \sigma(z_i)^{y_i} (1 - \sigma(z_i))^{1-y_i} \cdot \frac{1}{(2\pi)^{N/2} |\boldsymbol{\Sigma}(\mathbf{x}, \boldsymbol{\theta})|^{1/2}} \cdot \exp\left(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu}(\mathbf{x}, \boldsymbol{\theta}))^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}, \boldsymbol{\theta})(\mathbf{z} - \boldsymbol{\mu}(\mathbf{x}, \boldsymbol{\theta}))\right) \quad (7)$$

We consider two ways of inference and learning in GCRFBC model: (i) GCRFBCb - with conditional probability distribution $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})$, in which variables \mathbf{z} are marginalized over, and (ii) GCRFBCnb - with conditional probability distribution $P(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}, \mu_{\mathbf{z}})$, in which variables \mathbf{z} are substituted by their expectations.

3.3 INFERENCE IN GCRFBCB MODEL

Prediction of discrete outputs \mathbf{y} for given features \mathbf{x} and parameters $\boldsymbol{\theta}$ is analytically intractable due to integration of the joint distribution $P(\mathbf{y}, \mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$ with respect to latent variables. However, due to conditional independence between nodes, it is possible to obtain $P(y_i = 1|\mathbf{x}, \boldsymbol{\theta})$.

$$P(y_i = 1|\mathbf{x}, \boldsymbol{\theta}) = \int_{\mathbf{z}} \sigma(z_i) P(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta}) d\mathbf{z} \quad (8)$$

where $\sigma(z_i)$ models $P(y_i|\mathbf{z})$. As a result of independence properties of the distribution, it holds $P(y_i = 1|\mathbf{z}) = P(y_i = 1|z_i)$, and it is possible to marginalize $P(\mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$ with respect to latent variables $\mathbf{z}' = (z_1, \dots, z_{i-1}, z_{i+1}, \dots, z_N)$:

$$P(y_i = 1|\mathbf{x}, \boldsymbol{\theta}) = \int_{z_i} \sigma(z_i) \left(\int_{\mathbf{z}'} P(\mathbf{z}', z_i|\mathbf{x}, \boldsymbol{\theta}) d\mathbf{z}' \right) dz_i \quad (9)$$

where $\int_{z'} P(z', z_i | \mathbf{x}, \boldsymbol{\theta}) dz'$ is normal distribution with mean $\mu = \mu_i$ and variance $\sigma_i^2 = \Sigma_{ii}$. Therefore, it holds:

$$P(y_i = 1 | \mathbf{x}, \boldsymbol{\theta}) = \int_{-\infty}^{+\infty} \sigma(z_i) \mathcal{N}(z_i | \mu_i, \sigma_i^2) dz_i \quad (10)$$

The evaluation of $P(y_i = 0 | \mathbf{x}, \boldsymbol{\theta})$ is straightforward: $P(y_i = 0 | \mathbf{x}, \boldsymbol{\theta}) = 1 - P(y_i = 1 | \mathbf{x}, \boldsymbol{\theta})$.

The one-dimensional integral is still analytically intractable, but can be effectively evaluated by one-dimensional numerical integration. The proposed inference approach can be effectively used in case of huge number of nodes, due to low computational cost of one-dimensional numerical integration.

3.4 INFERENCE IN GCRFBCNB MODEL

The inference procedure in GCRFBCnb is much simpler, because marginalization with respect to latent variables is not performed. To predict \mathbf{y} , it is necessary to evaluate posterior maximum of latent variable $\mathbf{z}_{\max} = \underset{\mathbf{z}}{\operatorname{argmax}} P(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta})$, which is straightforward due to normal form of GCRF.

Therefore, it holds $\mathbf{z}_{\max} = \boldsymbol{\mu}_{\mathbf{z}, i}$. The conditional distribution $P(y_i = 1 | \mathbf{x}, \boldsymbol{\mu}_{\mathbf{z}, i}, \boldsymbol{\theta})$, where $\mu_{z,i}$ is expectation of latent variable z_i , can be expressed as:

$$P(y_i = 1 | \mathbf{x}, \boldsymbol{\mu}_{\mathbf{z}, i}, \boldsymbol{\theta}) = \sigma(\mu_{z,i}) = \frac{1}{1 + \exp(-\mu_{z,i})} \quad (11)$$

3.5 LEARNING IN GCRFBCB MODEL

In comparison with inference, learning procedure is more complicated. Evaluation of the conditional log likelihood is intractable, since latent variables cannot be analytically marginalized. The conditional log likelihood is expressed as:

$$\mathcal{L}(\mathbf{Y} | \mathbf{X}, \boldsymbol{\theta}) = \log \int_{\mathbf{Z}} P(\mathbf{Y}, \mathbf{Z} | \mathbf{X}, \boldsymbol{\theta}) d\mathbf{Z} = \sum_{j=1}^M \log \int_{z_j} P(\mathbf{y}_j, z_j | \mathbf{x}_j, \boldsymbol{\theta}) dz_j = \sum_{j=1}^M \mathcal{L}_j(\mathbf{y}_j | \mathbf{x}_j, \boldsymbol{\theta}) \quad (12)$$

$$\mathcal{L}_j(\mathbf{y}_j | \mathbf{x}_j, \boldsymbol{\theta}) = \log \int_{z_j} \prod_{i=1}^N \sigma(z_{ji})^{y_{ji}} (1 - \sigma(z_{ji}))^{1-y_{ji}} \frac{\exp(-\frac{1}{2}(\mathbf{z}_j - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{z}_j - \boldsymbol{\mu}_j))}{(2\pi)^{N/2} |\Sigma_j|^{1/2}} dz_j \quad (13)$$

where $\mathbf{Y} \in \mathbb{R}^{M \times N}$ is complete dataset of outputs, $\mathbf{X} \in \mathbb{R}^{M \times N \times A}$ is complete dataset of features, M is the total number of instances and A is the total number of features. Please note that each instance is structured, so while different instances are independent of each other, variables within one instance are dependent.

One way to approximate integral in conditional log likelihood is by local variational approximation. Jaakkola & Jordan (2000) derived lower bound for sigmoid function, which can be expressed as:

$$\sigma(x) \geq \sigma(\xi) \exp\{(x - \xi)/2 - \lambda(\xi)(x^2 - \xi^2)\} \quad (14)$$

where $\lambda(\xi) = -\frac{1}{2\xi} \cdot [\sigma(\xi) - \frac{1}{2}]$ and ξ is a variational parameter. The Eq. 14 is called ξ *transformation* of sigmoid function and it yields maximum value when $\xi = x$. This approximation can be applied to the model defined by Eq. 13, but the variational approximation has to be further extended because of the product of sigmoid functions, such that:

$$P(\mathbf{y}_j, z_j | \mathbf{x}_j, \boldsymbol{\theta}) = P(\mathbf{y}_j | z_j) P(z_j | \mathbf{x}_j, \boldsymbol{\theta}) \geq \underline{P}(\mathbf{y}_j, z_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) \quad (15)$$

$$\underline{P}(\mathbf{y}_j, z_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) = \prod_{i=1}^N \sigma(\xi_{ji}) \exp\left(z_{ji} y_{ji} - \frac{z_{ji} + \xi_{ji}}{2} - \lambda(\xi_{ji})(z_{ji}^2 - \xi_{ji}^2)\right) \cdot \frac{1}{(2\pi)^{N/2} |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{z}_j - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{z}_j - \boldsymbol{\mu}_j)\right) \quad (16)$$

The Eq. 16 can be arranged in the form suitable for integration. Detailed derivation of lower bound of conditional log likelihood is presented in Appendix A. The lower bound of conditional log likelihood

$\underline{\mathcal{L}}(y_j|x_j, \theta, \xi_j)$ is defined as:

$$\begin{aligned} \underline{\mathcal{L}}(y_j|x_j, \theta, \xi_j) = \log P(y_j|x_j, \theta, \xi_j) = \sum_{i=1}^N \left(\log \sigma(\xi_{ji}) - \frac{\xi_{ji}}{2} + \lambda(\xi_{ji})\xi_{ji}^2 \right) - \\ \frac{1}{2} \boldsymbol{\mu}_j^T \Sigma_j^{-1} \boldsymbol{\mu}_j + \frac{1}{2} \mathbf{m}_j^T S_j^{-1} \mathbf{m}_j + \frac{1}{2} \log |S_j| \end{aligned} \quad (17)$$

where:

$$S_j^{-1} = \Sigma_j^{-1} + 2\Lambda_j \quad \mathbf{m}_j = \Sigma_j \left(\left(y_j - \frac{1}{2} \mathbf{1} \right) + \Sigma_j^{-1} \boldsymbol{\mu}_j \right) \quad (18)$$

$$\Lambda_j = \begin{bmatrix} \lambda(\xi_{j1}) & 0 & 0 & \dots & 0 \\ 0 & \lambda(\xi_{j2}) & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \lambda(\xi_{jN}) \end{bmatrix} \quad (19)$$

GCRFBCb uses the derivative of conditional log likelihood in order to find the optimal values for parameters α , β and matrix of variational parameters $\boldsymbol{\xi} \in \mathbb{R}^{M \times N}$. In order to ensure positive definiteness of normal distribution involved, it is sufficient to constrain parameters $\alpha > 0$ and $\beta > 0$. The partial derivatives of lower bound of conditional log likelihood are presented in Appendix B. For constrained optimization, the truncated Newton algorithm was used Nocedal & Wright (2006); Facchinei et al. (2002). The target function is not convex, so finding a global optimum cannot be guaranteed.

3.6 LEARNING IN GCRFBCNB MODEL

In GCRFBCnb the mode of posterior distribution of continuous latent variable \mathbf{z} is evaluated directly, so there is no need for approximation. The conditional log likelihood can be expressed as:

$$\mathcal{L}(\mathbf{Y}|\mathbf{X}, \theta, \boldsymbol{\mu}) = \log P(\mathbf{Y}|\mathbf{X}, \theta, \boldsymbol{\mu}) = \sum_{j=1}^M \sum_{i=1}^N \log P(y_{ji}|x_j, \theta, \mu_{ji}) = \sum_{j=1}^M \sum_{i=1}^N \underline{\mathcal{L}}_{ji}(y_{ji}|x_j, \theta, \mu_{ji}) \quad (20)$$

$$\underline{\mathcal{L}}_{ji}(y_{ji}|x_j, \theta, \mu_{ji}) = y_{ji} \log \sigma(\mu_{ji}) + (1 - y_{ji}) \log (1 - \sigma(\mu_{ji})) \quad (21)$$

The partial derivatives of conditional log likelihood are presented in Appendix C.

4 EXPERIMENTAL EVALUATION

Both proposed models were tested and compared on synthetic data and real-world tasks.¹ All compared classifiers were compared in terms of the area under ROC curve (AUC) and accuracy² (ACC). Moreover, the lower bound (in case of GCRFBCb) of conditional log likelihood $\underline{\mathcal{L}}(\mathbf{Y}|\mathbf{X}, \theta, \boldsymbol{\mu})$ and actual value (in case of GCRFBCnb) of conditional log likelihood $\mathcal{L}(\mathbf{Y}|\mathbf{X}, \theta)$ of obtained values on synthetic test dataset were also reported.

4.1 SYNTHETIC DATASET

The main goal of experiments on synthetic datasets was to examine models under various controlled conditions, and show advantages and disadvantages of each. In all experiments on synthetic datasets two different graphs were used (hence $\beta \in \mathbb{R}^2$) and two unstructured predictors (hence $\alpha \in \mathbb{R}^2$). The results of experiments on synthetic datasets are presented in Appendix D.

It can be noticed, that in cases where norm of the variances of latent variables is small, both models have equal performance considering AUC and conditional log likelihood $\underline{\mathcal{L}}(\mathbf{Y}|\mathbf{X}, \theta)$. This is the case when values of parameters α used in data generating process are greater or equal to the

¹Implementation can be found at https://github.com/andrijaster/GCRFBC_B_NB

²PyStruct package does not have option of returning SSVM and CRF confidence values for AUC evaluation

values of parameters β . This means that the information provided by unstructured predictors is more important for classifications task than the information provided by output structure. Therefore, conditional distribution $P(\mathbf{y}, \mathbf{z}|\mathbf{x}, \boldsymbol{\theta})$ is concentrated around mean value and MAP estimate is a satisfactory approximation. However, when data is generated from distribution with significantly higher values of β than α , the GCRFBCb performs significantly better than GCRFBCnb. For the larger values of variance norm, this difference is also large. This means that the structure between outputs has significant contribution to solving the classification task. It can be concluded that GCRFBCb has at least equal prediction performance as GCRFBCnb. Also, it can be argued that the models were generally able to utilize most of the information (from both features and the structure between outputs), which can be seen through AUC values. In addition, distribution of local variational parameters were analyzed during learning. It is noticed that in each epoch, the variance of this distribution is small and that the parameters can be clustered and their number significantly reduced. Therefore, it is possible to significantly lower down computational and memory costs of GCRFBCb learning procedure, but that's out of the scope of this paper.

4.2 PERFORMANCE ON REAL-WORLD DATASETS

4.2.1 SKI LIFTS CONGESTION

Data used in this research includes information on ski lift gate entrances in Kopaonik ski resorts, for the period March 15 to March 30 for the seasons from 2006 to 2011. The goal is to predict occurrence of crowding on ski lifts for 40 minutes in advance. Total number of instances in dataset was 4,850 for each ski lift, which is 33,950 in total.

Relatively simple method for crowding detection was devised for labelling data. We assume that, if the crowding at some gate occurs, distributions of skiing times from other gates to that gate within some time window get shifted towards larger values. We model probability distribution of skiing time between two gates by the well-known parametric method of kernel density estimation (KDE) (Silverman, 2018). The distribution shift is measured with respect to the mode of the distribution. The dataset is generated by observing shifts in time windows of 5 minutes. When the mode of the distribution of skiing times within that window is greater than the mode for the whole time-span, the instance is labeled by 1 (crowding) and otherwise, it is labeled by 0 (no crowding). In order to obtain more information from the data distribution, additional 18 features were extracted.

Four different unstructured predictors that were trained on each class separately were used: ridge logistic regression, LASSO logistic regression, neural network and random forest, whereas additional two unstructured predictors: decision tree and neural network were trained on all nodes together. Additionally, three structural support vector machine and two CRFs classifiers were used (Müller & Behnke, 2014). Fully connected graph of SSVM and CRF models are defined as SSVM-full and CRF-full, whereas Chow-Liu tree method for specifying edge connections are defined as SSVM-tree and CRF-tree, respectively. In the SSVM-independent model the nodes of the graph are not connected.

Six different weighted graphs were used to capture dependence structure between ski lifts (nodes): χ^2 statistics on labels of training set, mutual information between labels, correlation matrix between outputs of over-fitted neural networks, norm of difference between vectors of labels and two graphs were defined based on difference of vectors of historical labels and on differences of historical averages of skier times.

The AUC score and ACC of structured and unstructured predictors, along with the total computational time are shown in Table 1. It can be observed that GCRFBCb and GCRBCnb outperformed unstructured and other structured predictors in all cases. Based on evaluated parameters it could be concluded that dependence structure has significant impact on overall prediction performance, even though, due to low values of norm of variance, GCRFBCb and GCRFBCnb have equal AUC scores. It can be summarized that advantages of structured models compared to unstructured are obvious, but in this particular task due to equal prediction performance and its lower computational and memory complexity, GCRFBCnb is the best choice for this specific application.

Table 1: Prediction performance and computation time of classifiers - Ski lifts congestion problem

Model	AUC	ACC	Calculation time [sec]
GCRFBCnb	0.831	0.749	119.554
GCRFBCb	0.831	0.749	3364.326
Ridge logistic	0.793	0.736	0.41
LASSO logistic	0.793	0.735	1.799
Neural network	0.790	0.720	151.571
Random forest	0.783	0.720	7.983
Decision tree - together	-	0.681	8.297
Neural network - together	-	0.711	13.997
SSVM - full	-	0.622	517.412
SSVM - tree	-	0.615	580.475
SSVM - independent	-	0.635	1029.172
CRF - tree	-	0.745	16415.723
CRF - full	-	0.740	13942.542

Table 2: Prediction performance and computation time of classifiers - Music classification according to emotion

Model	AUC	ACC	Calculation time [sec]
GCRFBCnb	0.859	0.811	7.248
GCRFBCb	0.860	0.813	353.328
Ridge logistic	0.826	0.794	0.138
LASSO logistic	0.832	0.797	0.874
Neural network	0.811	0.783	98.132
Random forest	0.843	0.798	2.469
Decision tree - together	-	0.736	0.564
Neural network - together	-	0.782	8.471
SSVM - full	-	0.755	76.817
SSVM - tree	-	0.795	75.93
SSVM - independent	-	0.784	146.867

4.2.2 MULTI-LABEL CLASSIFICATION OF MUSIC ACCORDING TO EMOTION

The dataset used for this work consists of 100 songs from 7 different genres. The collection was created from 233 musical albums choosing three songs from each album. 8 rhythmic and 64 timbre features are extracted. The music is labeled in 6 categories of emotions: amazed-surprised, happy-pleased, relaxing-calm, quiet-still, sad-lonely and angry-fearful (Trohidis et al., 2008). Total number of instances in dataset was 593. Four different weighted graphs were used: statistics on labels of training set, mutual information between labels, correlation matrix between outputs of over-fitted neural networks and norm of difference between vectors of labels. Same unstructured predictors as in ski lift congestion problem were used, along with three structural support vector machine classifiers.

The performances of models are evaluated by 10 fold cross validation. The AUC score and ACC of structured and unstructured predictors, along with the total computational time are shown in Table 2. It can be seen that GCRFBCb has achieved the best prediction performances. The ACC of GCRFBC models are significantly better than the SSVM performances. The AUC score and ACC of GCRFBCb are higher than the best result (AUC = 0.8237) presented in original paper (Trohidis et al., 2008). As in previous cases, computational time of GCRFBCb is significantly longer compared to GCRFBCnb and SSVM models.

4.2.3 GENE FUNCTION CLASSIFICATION

This dataset is formed by micro-array expression data and phylogenetic profiles with 2417 genes (instances). The number of features is 103, whereas each gene is associated with the set of 14 groups (Elisseeff & Weston, 2002). The same unstructured, structured predictors and weighted graphs, as

Table 3: Prediction performance and computation time of classifiers - Gene classification problem

Model	AUC	ACC	Calculation time [sec]
GCRFBCnb	0.775	0.766	48.167
GCRFBCb	0.797	0.775	2297.727
Ridge logistic	0.582	0.539	0.079
LASSO logistic	0.583	0.540	0.188
Neural network	0.580	0.567	70.298
Random forest	0.601	0.615	5.529
Decision tree - together	-	0.691	1.218
Neural network - together	-	0.775	28.381
SSVM - full	-	0.771	10137.049
SSVM - tree	-	0.768	722.156
SSVM - independent	-	0.539	78.8870

in music according to emotion classification, were used. The 10-fold cross validation results of the classification are shown in Table 3.

It can be observed that both GCRFBCb and GCRFBCnb achieved significantly better results in comparison with unstructured predictors. However, neural network trained on all data together achieved the same ACC scores as GCRFBCb. The AUC of GCRFBCb has outperformed Random forest classifier by 19%, whereas SSVM - tree has better ACC compared to GCRFBCnb. It also outperformed GCRFBCnb, but as expected, its computation time was longer. In addition, the computation time of CRFs models are longer compared to GCRFBCb

4.2.4 HIGHWAY CONGESTION

The E70-E75 motorway is a major transit motorway in Serbia. With 504 kilometers, it is the one the major transit motorway in Serbia. It crosses the country from north-west to south, starting at Batrovci border crossing with the Republic of Croatia and ending with Preševo border crossing with the Republic of North Macedonia.

One of the biggest problems in E70-E75 motorway is high congestion that frequently occurs. One of the reasons lies in lack of open toll stations. In order to mitigate congestion problem, it is necessary to predict its occurrence and open enough toll stations. Data used in this research includes information of car entrance and exit for the year 2017. Two different sections were analyzed: Belgrade - Adaševci and Niš - Belgrade. The section Belgrade - Adaševci was analyzed for the period of January 2017, whereas section Niš - Belgrade was analyzed for the period of April - July 2017. The congestion was labeled using the similar technique based on KDE as presented in the ski lifts congestion problem. Based on raw datasets for sections Niš - Belgrade and Belgrade - Adaševci with 5,132,918 and 487,767 instances, respectively, a new dataset for section Niš - Belgrade is generated by observing shifts in time windows of 10 minutes due to large number of vehicles, whereas in the case of section Belgrade - Adaševci the shifts are observed in time windows of 20 minutes. Total numbers of instances for sections Belgrade - Adaševci and Niš - Belgrade are 50,964 and 235,872, whereas numbers of highway exits (outputs) are 6 and 18, respectively. The extracted features are similar to the ones presented in ski congestion problem. The χ^2 statistics, mutual information, correlation matrix and difference of vectors of historical labels were used to capture dependence structure, whereas the same unstructured predictors as in ski lifts congestion problem were evaluated. The classification results, validated by 10 fold cross validation, are presented in Table 4.

The GCRFBCnb achieved the highest AUC and ACC scores in the section Belgrade - Adaševci, whereas GCRFBCb has better prediction performance in section Niš - Belgrade. Moreover, in case of section Niš - Belgrade, GCRFBCb has worse ACC score than fully connected CRF, whereas CRF-tree outperformed GCRFBCnb in section Belgrade - Adaševci

Table 4: Prediction performance and computation time of classifiers - Highway congestion problem

	Niš - Belgrade			Belgrade - Adaševci		
	AUC	ACC	Calculation time [sec]	AUC	ACC	Calculation time [sec]
GCRFBCnb	0.740	0.684	344.166	0.974	0.925	90.321
GCRFBCb	0.751	0.692	13818.874	0.956	0.895	2103.749
Ridge logistic	0.716	0.681	10.73	0.917	0.856	1.771
LASSO logistic	0.716	0.680	30.12	0.917	0.856	1.657
Neural network	0.72	0.682	857.602	0.956	0.904	125.339
Random forest	0.739	0.683	209.589	0.965	0.914	3.826
Decision tree - together	-	0.625	635.464	-	0.898	1.893
Neural network - together	-	0.664	125.441	-	0.880	16.475
SSVM - full	-	0.588	7637.794	-	0.739	340.806
SSVM - tree	-	0.588	3684.138	-	0.755	392.597
SSVM - independent	-	0.602	3262.208	-	0.814	704.07
CRF - tree	-	0.685	29749.054	-	0.88	26539.250
CRF - full	-	0.683	52563.972	-	0.898	25339.97

5 CONCLUSION

In this paper, a new model, called Gaussian Conditional Random Fields for Binary Classification (GCRFBC) is presented. The model is based on latent GCRF structure, which means that intractable structured classification problem can become tractable and efficiently solved. Moreover, the improvements previously applied to regression GCRF can be easily extended to GCRFBC. Two different variants of GCRFBC were derived: GCRFBCb and GCRFBCnb. Empirical Bayes (marginalization of latent variables) by local variational methods is used in optimization procedure of GCRFBCb, whereas MAP estimate of latent variables is applied in GCRFBCnb. Based on presented methodology and obtained experimental results on synthetic and real-world datasets it can be concluded that both GCRFBCb and GCRFBCnb models have better prediction performance compared to the analysed structured unstructured predictors. Additionally, GCRFBCb has better performance considering AUC score, ACC and lower bound of conditional log likelihood $\mathcal{L}(Y|X, \theta)$ compared to GCRFBCnb, in cases where norm of the variances of latent variables is high. However, in cases where norm of the variances is close to zero, both models have equal prediction performance. Due to high memory and computational complexity of GCRFBCb compared to GCRFBCnb, in cases where norm of the variances is close to zero, it is reasonable to use GCRFBCnb. Additionally, the trade off between complexity and accuracy can be made in situation where norm of the variances is high. Further studies should address extending GCRFBC to structured multi-label classification problems, and lower computational complexity of GCRFBCb by considering efficient approximations.

REFERENCES

- David Belanger and Andrew McCallum. Structured prediction energy networks. In *International Conference on Machine Learning*, pp. 983–992, 2016.
- David Belanger, Bishan Yang, and Andrew McCallum. End-to-end learning for structured prediction energy networks. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pp. 429–439. JMLR. org, 2017.
- Gang Chen, Yawei Li, and Sargur N Srihari. Word recognition with deep conditional random fields. *arXiv preprint arXiv:1612.01072*, 2016.
- Ryan Cotterell and Kevin Duh. Low-resource named entity recognition with cross-lingual, character-level neural conditional random fields. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, volume 2, pp. 91–96, 2017.
- André Elisseeff and Jason Weston. A kernel method for multi-labelled classification. In *Advances in neural information processing systems*, pp. 681–687, 2002.
- Francisco Facchinei, Stefano Lucidi, and Laura Palagi. A truncated newton algorithm for large scale box constrained optimization. *SIAM Journal on Optimization*, 12(4):1100–1125, 2002.

- John Fox. *Applied regression analysis and generalized linear models*. Sage Publications, 2015.
- Benjamin Frot, Luke Jostins, and Gilean McVean. Graphical model selection for gaussian conditional random fields in the presence of latent variables. *Journal of the American Statistical Association*, (just-accepted), 2018.
- Jesse Glass, Mohamed F Ghalwash, Milan Vukicevic, and Zoran Obradovic. Extending the modelling capacity of gaussian conditional random fields while learning faster. In *AAAI*, pp. 1596–1602, 2016.
- Tommi S Jaakkola and Michael I Jordan. Bayesian parameter estimation via variational methods. *Statistics and Computing*, 10(1):25–37, 2000.
- Minyoung Kim. Mixtures of conditional random fields for improved structured output prediction. *IEEE transactions on neural networks and learning systems*, 28(5):1233–1240, 2017.
- Sergey Kosov, Kimiaki Shirahama, Chen Li, and Marcin Grzegorzec. Environmental microorganism classification using conditional random fields and deep convolutional neural networks. *Pattern Recognition*, 77:248–261, 2018.
- You Lu and Bert Huang. Structured output learning with conditional generative flows. *arXiv preprint arXiv:1905.13288*, 2019.
- Laurens Maaten, Max Welling, and Lawrence Saul. Hidden-unit conditional random fields. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 479–488, 2011.
- Kristen Masada and Razvan C Bunescu. Chord recognition in symbolic music using semi-markov conditional random fields. In *ISMIR*, pp. 272–278, 2017.
- Andreas C Müller and Sven Behnke. Pystruct: learning structured prediction in python. *The Journal of Machine Learning Research*, 15(1):2055–2060, 2014.
- Jorge Nocedal and Stephen J Wright. *Numerical optimization* 2nd, 2006.
- Vladan Radosavljevic. *Gaussian conditional random fields for regression in remote sensing*. Temple University, 2011.
- Vladan Radosavljevic, Slobodan Vucetic, and Zoran Obradovic. Continuous conditional random fields for regression in remote sensing. In *ECAI*, pp. 809–814, 2010.
- Vladan Radosavljevic, Slobodan Vucetic, and Zoran Obradovic. Neural gaussian conditional random fields. In *Joint European conference on machine learning and knowledge discovery in databases*, pp. 614–629. Springer, 2014.
- Bernard W Silverman. *Density estimation for statistics and data analysis*. Routledge, 2018.
- Jelena Stojanovic, Milos Jovanovic, Djordje Gligorijevic, and Zoran Obradovic. Semi-supervised learning for structured regression on partially observed attributed graphs. In *Proceedings of the 2015 SIAM International Conference on Data Mining*, pp. 217–225. SIAM, 2015.
- Gilbert Strang, Gilbert Strang, Gilbert Strang, and Gilbert Strang. *Introduction to linear algebra*, volume 3. Wellesley-Cambridge Press Wellesley, MA, 1993.
- Hongyu Su. *Multilabel Classification through Structured Output Learning - Methods and Applications*. Aalto University, 2015.
- Charles Sutton and Andrew McCallum. *An introduction to conditional random fields for relational learning*, volume 2. Introduction to statistical relational learning. MIT Press, 2006.
- Chenhao Tan, Jie Tang, Jimeng Sun, Quan Lin, and Fengjiao Wang. Social action tracking via noise tolerant time-varying factor graphs. In *Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 1049–1058. ACM, 2010.

- Marshall F Tappen, Ce Liu, Edward H Adelson, and William T Freeman. Learning gaussian conditional random fields for low-level vision. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1–8. IEEE, 2007.
- Konstantinos Trohidis, Grigorios Tsoumakas, George Kalliris, and Ioannis P Vlahavas. Multi-label classification of music into emotions. In *ISMIR*, volume 8, pp. 325–330, 2008.
- Matt Wytock and Zico Kolter. Sparse gaussian conditional random fields: Algorithms, theory, and application to energy forecasting. In *International conference on machine learning*, pp. 1265–1273, 2013.
- Peng Zhang, Ming Li, Yan Wu, and Hejing Li. Hierarchical conditional random fields model for semisupervised sar image segmentation. *IEEE Transactions on Geoscience and Remote Sensing*, 53(9):4933–4951, 2015.
- Haris Bin Zia, Agha Ali Raza, and Awais Athar. Urdu word segmentation using conditional random fields (crfs). *arXiv preprint arXiv:1806.05432*, 2018.

A DERIVATION OF LOWER BOUND OF CONDITIONAL LIKELIHOOD

In this section we derive lower bound of conditional likelihood. In order to obtain suitable form of joint distribution that can be easily integrated, the lower bound for sigmoid function was used (Jaakkola & Jordan, 2000). The lower bound of joint distribution $P(\mathbf{y}_j, \mathbf{z}_j | \mathbf{x}_j, \boldsymbol{\theta})$ can be expressed as:

$$P(\mathbf{y}_j, \mathbf{z}_j | \mathbf{x}_j, \boldsymbol{\theta}) = P(\mathbf{y}_j | \mathbf{z}_j) P(\mathbf{z}_j | \mathbf{x}_j, \boldsymbol{\theta}) \geq \underline{P}(\mathbf{y}_j, \mathbf{z}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) \quad (22)$$

$$\underline{P}(\mathbf{y}_j, \mathbf{z}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) = \prod_{i=1}^N \sigma(\xi_{ji}) \exp \left(z_{ji} y_{ji} - \frac{z_{ji} + \xi_{ji}}{2} - \lambda(\xi_{ji})(z_{ji}^2 - \xi_{ji}^2) \right) \cdot \frac{1}{(2\pi)^{N/2} |\Sigma_j|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{z}_j - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} (\mathbf{z}_j - \boldsymbol{\mu}_j) \right) \quad (23)$$

The simplified form of Eq. 23 can be represented by rearranging terms in the following form:

$$\underline{P}(\mathbf{y}_j, \mathbf{z}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) = \mathcal{T}(\boldsymbol{\xi}_j) \exp \left(\mathbf{z}_j^T (\mathbf{y}_j - \frac{1}{2} \mathbf{I}) - \lambda \mathbf{z}_j^T \mathbf{z}_j - \frac{1}{2} \mathbf{z}_j^T \Sigma_j^{-1} \mathbf{z}_j + \mathbf{z}_j^T \Sigma_j^{-1} \boldsymbol{\mu}_j \right) \quad (24)$$

$$\mathcal{T}(\boldsymbol{\xi}_j) = \frac{1}{(2\pi)^{N/2} |\Sigma_j|^{1/2}} \prod_{i=1}^N \sigma(\xi_{ji}) \exp \left(-\frac{1}{2} \boldsymbol{\mu}_j^T \Sigma_j^{-1} \boldsymbol{\mu}_j - \frac{\xi_{ji}}{2} + \lambda(\xi_{ji}) \xi_{ji}^2 \right) \quad (25)$$

The lower bound of likelihood $\underline{P}(\mathbf{y}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j)$ can be obtained by marginalization of \mathbf{z}_j as:

$$\begin{aligned} \underline{P}(\mathbf{y}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) &= \int \underline{P}(\mathbf{y}_j, \mathbf{z}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) d\mathbf{z}_j \\ &= \mathcal{T}(\boldsymbol{\xi}_j) \int \exp \left(\mathbf{z}_j^T (\mathbf{y}_j - \frac{1}{2} \mathbf{I}) - \lambda \mathbf{z}_j^T \mathbf{z}_j - \frac{1}{2} \mathbf{z}_j^T \Sigma_j^{-1} \mathbf{z}_j + \mathbf{z}_j^T \Sigma_j^{-1} \boldsymbol{\mu}_j \right) d\mathbf{z}_j \\ &= \mathcal{T}(\boldsymbol{\xi}_j) \int \exp \left(-\frac{1}{2} \mathbf{z}_j^T (\Sigma_j^{-1} + 2\Lambda_j) \mathbf{z}_j + \mathbf{z}_j^T (\Sigma_j^{-1} + 2\Lambda_j) (\mathbf{y}_j - \frac{1}{2} \mathbf{I}) + \mathbf{z}_j^T \Sigma_j^{-1} \boldsymbol{\mu}_j \right) d\mathbf{z}_j \end{aligned} \quad (26)$$

The lower bound of likelihood $\underline{P}(\mathbf{y}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j)$ can be transformed in the following form:

$$\begin{aligned} \underline{P}(\mathbf{y}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) &= \mathcal{T}(\boldsymbol{\xi}_j) \int \exp \left(-\frac{1}{2} (\mathbf{z}_j - \mathbf{m}_j)^T S_j^{-1} (\mathbf{z}_j - \mathbf{m}_j) + \frac{1}{2} \mathbf{m}_j^T S_j^{-1} \mathbf{m}_j \right) d\mathbf{z}_j \\ &= \mathcal{T}(\boldsymbol{\xi}_j) \exp \left(\frac{1}{2} \mathbf{m}_j^T S_j^{-1} \mathbf{m}_j \right) \int \exp \left(-\frac{1}{2} (\mathbf{z}_j - \mathbf{m}_j)^T S_j^{-1} (\mathbf{z}_j - \mathbf{m}_j) \right) d\mathbf{z}_j \end{aligned} \quad (27)$$

where $S_j^{-1} = \Sigma_j^{-1} + 2\Lambda_j$ and $\mathbf{m}_j = \Sigma_j \left((\mathbf{y}_j - \frac{1}{2} \mathbf{I}) + \Sigma_j^{-1} \boldsymbol{\mu}_j \right)$.

This integration is easily performed by noting that it is the integral over an unnormalized Gaussian distribution, which yields:

$$\underline{P}(\mathbf{y}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) = (2\pi)^{N/2} |\Sigma_j|^{1/2} \mathcal{T}(\boldsymbol{\xi}_j) \exp \left(\frac{1}{2} \mathbf{m}_j^T S_j^{-1} \mathbf{m}_j \right) |\Sigma_j|^{1/2} \quad (28)$$

The final form of the lower bound of conditional log likelihood $\underline{\mathcal{L}}_j(\mathbf{y}_j | \mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j)$ is:

$$\begin{aligned} \underline{\mathcal{L}}_j(\mathbf{y}_j|\mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) &= \log \underline{P}(\mathbf{y}_j|\mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j) = \sum_{i=1}^N \left(\log \sigma(\xi_{ji}) - \frac{\xi_{ji}}{2} + \lambda(\xi_{ji})\xi_{ji}^2 \right) - \\ &\quad \frac{1}{2} \boldsymbol{\mu}_j^T \Sigma_j^{-1} \boldsymbol{\mu}_j + \frac{1}{2} \mathbf{m}_j^T S_j^{-1} \mathbf{m}_j + \frac{1}{2} \log |S_j| \end{aligned} \quad (29)$$

B PARTIAL DERIVATIVE OF LOWER BOUND OF CONDITIONAL LOG LIKELIHOOD

The partial derivative of lower bound of conditional log likelihood (GCRFBCb) $\frac{\partial \underline{\mathcal{L}}_j(\mathbf{y}_j|\mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j)}{\partial \alpha_k}$ is computed as:

$$\begin{aligned} \frac{\partial \underline{\mathcal{L}}_j(\mathbf{y}_j|\mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j)}{\partial \alpha_k} &= -\frac{1}{2} \text{Tr} \left(S_j \frac{\partial S_j^{-1}}{\partial \alpha_k} \right) + \frac{\partial \mathbf{m}_j^T}{\partial \alpha_k} S_j^{-1} \mathbf{m}_j + \frac{1}{2} \mathbf{m}_j^T \frac{\partial S_j^{-1}}{\partial \alpha_k} \mathbf{m}_j \\ &\quad - \frac{\boldsymbol{\mu}_j^T}{\partial \alpha_k} \Sigma_j^{-1} \boldsymbol{\mu}_j - \frac{1}{2} \boldsymbol{\mu}_j^T \frac{\partial \Sigma_j^{-1}}{\partial \alpha_k} + \frac{1}{2} \text{Tr} \left(\Sigma_j \frac{\partial \Sigma_j^{-1}}{\partial \alpha_k} \right) \end{aligned} \quad (30)$$

where:

$$\frac{\partial S_j^{-1}}{\partial \alpha_k} = \frac{\partial \Sigma_j^{-1}}{\partial \alpha_k} = \begin{cases} 2, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases} \quad (31)$$

$$\frac{\partial \mathbf{m}_j^T}{\partial \alpha_k} = - \left(\mathbf{y}_j - \frac{1}{2} \mathbf{I} + \boldsymbol{\mu}_j^T \Sigma_j^{-1} \right) S_j \frac{\partial S_j^{-1}}{\partial \alpha_k} S_j + \frac{\partial \boldsymbol{\mu}_j^T}{\partial \alpha_k} \Sigma_j^{-1} S_j + \boldsymbol{\mu}_j^T \frac{\partial \Sigma_j^{-1}}{\partial \alpha_k} S_j \quad (32)$$

$$\frac{\partial \boldsymbol{\mu}_j^T}{\partial \alpha_k} = \left(2\alpha_k R_k(\mathbf{x}) - \frac{\partial \Sigma_j^{-1}}{\partial \alpha_k} \boldsymbol{\mu}_j \right)^T \Sigma_j^T \quad (33)$$

Similarly partial derivatives with respect to β can be defined as:

$$\begin{aligned} \frac{\partial \underline{\mathcal{L}}_j(\mathbf{y}_j|\mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j)}{\partial \beta_l} &= -\frac{1}{2} \text{Tr} \left(S_j \frac{\partial S_j^{-1}}{\partial \beta_l} \right) + \frac{\partial \mathbf{m}_j^T}{\partial \beta_l} S_j^{-1} \mathbf{m}_j + \frac{1}{2} \mathbf{m}_j^T \frac{\partial S_j^{-1}}{\partial \beta_l} \mathbf{m}_j \\ &\quad - \frac{\boldsymbol{\mu}_j^T}{\partial \beta_l} \Sigma_j^{-1} \boldsymbol{\mu}_j - \frac{1}{2} \boldsymbol{\mu}_j^T \frac{\partial \Sigma_j^{-1}}{\partial \beta_l} + \frac{1}{2} \text{Tr} \left(\Sigma_j \frac{\partial \Sigma_j^{-1}}{\partial \beta_l} \right) \end{aligned} \quad (34)$$

where:

$$\frac{\partial S_j^{-1}}{\partial \beta_l} = \frac{\partial \Sigma_j^{-1}}{\partial \beta_l} = \begin{cases} \sum_{n=1}^N e_{in}^l S_{in}^l(x), & \text{if } i = j \\ -e_{ij}^l S_{ij}^l(x), & \text{if } i \neq j \end{cases} \quad (35)$$

$$\frac{\partial \mathbf{m}_j^T}{\partial \beta_l} = - \left(\mathbf{y}_j - \frac{1}{2} \mathbf{I} + \boldsymbol{\mu}_j^T \Sigma_j^{-1} \right) S_j \frac{\partial S_j^{-1}}{\partial \beta_l} S_j + \frac{\partial \boldsymbol{\mu}_j^T}{\partial \beta_l} \Sigma_j^{-1} S_j + \boldsymbol{\mu}_j^T \frac{\partial \Sigma_j^{-1}}{\partial \beta_l} S_j \quad (36)$$

$$\frac{\partial \boldsymbol{\mu}_j^T}{\partial \beta_l} = \left(-\frac{\partial \Sigma_j^{-1}}{\partial \beta_l} \boldsymbol{\mu}_j \right)^T \Sigma_j^T \quad (37)$$

In the same manner partial derivatives of conditional log likelihood with respect to ξ_{ji} are:

$$\begin{aligned} \frac{\partial \underline{\mathcal{L}}_j(\mathbf{y}_j|\mathbf{x}_j, \boldsymbol{\theta}, \boldsymbol{\xi}_j)}{\partial \xi_{ji}} &= -\frac{1}{2} \text{Tr} \left(2S_j \frac{\partial \Lambda_j}{\partial \xi_{ji}} \right) - \left[2 \left(\mathbf{y}_j - \frac{1}{2} \mathbf{I} \right) S_j \frac{\partial \Lambda_j}{\partial \xi_{ji}} S_j \right] S_j^{-1} \mathbf{m}_j \\ &\quad + \mathbf{m}_j^T \frac{\partial \Lambda_j}{\partial \xi_{ji}} \mathbf{m}_j + \sum_{i=1}^N \left(\left(\frac{1}{\sigma(\xi_{ji})} + \frac{1}{2} \xi_{ji} \right) \frac{\partial \sigma(\xi_{ji})}{\partial \xi_{ji}} + \frac{1}{2} \left(\sigma(\xi_{ji}) - \frac{3}{4} \right) \right) \end{aligned} \quad (38)$$

where:

$$\frac{\partial \Lambda_j}{\partial \xi_{ji}} = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ 0 & 0 & \frac{\partial \lambda(\xi_{ji})}{\partial \xi_{ji}} & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} \quad (39)$$

$$\frac{\partial \sigma(\xi_{ji})}{\partial \xi_{ji}} = \sigma(\xi_{ji})(1 - \sigma(\xi_{ji})) \quad (40)$$

$$\frac{\partial \lambda(\xi_{ji})}{\partial \xi_{ji}} = \frac{1}{2\xi_{ji}} \frac{\partial \sigma(\xi_{ji})}{\partial \xi_{ji}} - \frac{1}{2} \left(\sigma(\xi_{ji}) - \frac{1}{2} \right) \frac{1}{\xi_{ji}^2} \quad (41)$$

C PARTIAL DERIVATIVE OF CONDITIONAL LOG LIKELIHOOD

The derivatives of the conditional log likelihood (GCRFBCnb) with respect to α and β are defined as, respectively:

$$\frac{\partial \mathcal{L}_{ji}(y_{ji} | \mathbf{x}_j, \boldsymbol{\theta}, \mu_{ji})}{\partial \alpha_k} = (y_{ji} - \sigma(\mu_{ji})) \frac{\partial \mu_{ji}}{\partial \alpha_k} \quad (42)$$

$$\frac{\partial \mathcal{L}_{ji}(y_{ji} | \mathbf{x}_j, \boldsymbol{\theta}, \mu_{ji})}{\partial \beta_l} = (y_{ji} - \sigma(\mu_{ji})) \frac{\partial \mu_{ji}}{\partial \beta_l} \quad (43)$$

where $\frac{\partial \mu_{ji}}{\partial \alpha_k}$ and $\frac{\partial \mu_{ji}}{\partial \beta_l}$ are elements of the vectors $\frac{\partial \boldsymbol{\mu}_j}{\partial \boldsymbol{\alpha}_k}$ and $\frac{\partial \boldsymbol{\mu}_j}{\partial \boldsymbol{\beta}_l}$ and can be obtained by Eqs. 33 and 37, respectively.

D SYNTHETIC DATASET RESULTS

In order to generate and label graph nodes, edge weights S and unstructured predictor values R were randomly generated from uniform distribution. Besides, it was necessary to choose values of parameters α and β . Greater values of α indicate that the model is more confident about performance of unstructured predictors, whereas for the larger value of β the model is putting more emphasis on the dependence structure of output variables.

Six different values of parameters α and β were used. In the first group α and β have similar values, so unstructured predictors and dependence structure between outputs have similar importance. In the second group, α has higher values compared to β , which means that unstructured predictors are more important than the dependence structure. In the third group β has higher values than α , meaning that dependence structure is more important than unstructured predictors.

Along with the AUC and conditional log likelihood, norm of the variances of latent variables (diagonal elements in the covariance matrix) is evaluated and presented in Table 5. In addition, the results of experiments are presented in Fig. 1, where for different values of α and β we show differences between GCRFBCb and GCRFBCnb (a) AUC scores, (b) log likelihoods, and (c) norm of the variances of latent variables.

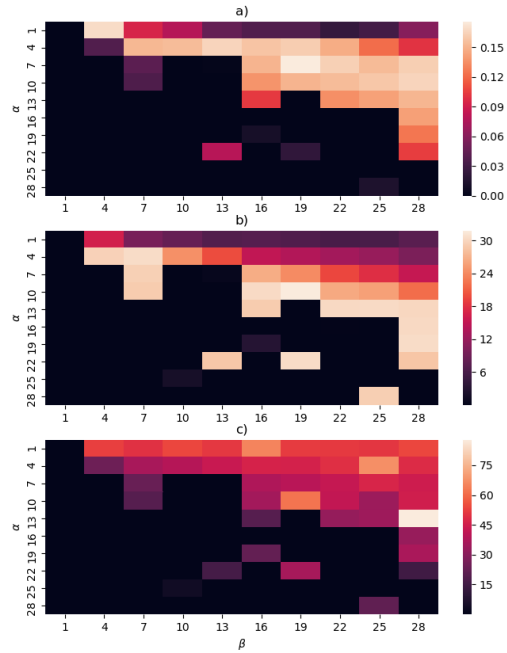


Figure 1: Experimental evaluation of differences between GCRFBCb and GCRFBCnb (a) AUC scores, (b) log likelihoods, and (c) norms of the variances of latent variables for different values of α and β

Table 5: Comparison of GCRFBCb and GCRFBCnb prediction performance for different values of α and β , as measured by AUC, log likelihood, and norm of diagonal elements of the covariance matrix

No.	Parameters	GCRFBCb			GCRFBCnb	
		AUC	$\mathcal{L}(\mathbf{Y} \mathbf{X}, \theta)$	$\ \sigma\ _2$	AUC	$\mathcal{L}(\mathbf{Y} \mathbf{X}, \theta)$
1	$\alpha = [5, 4]$ $\beta = [5, 22]$	0.812	-71.150	0.000	0.812	-71.151
2	$\alpha = [1, 18]$ $\beta = [1, 18]$	0.903	-75.033	0.001	0.902	-75.033
3	$\alpha = [22, 21]$ $\beta = [5, 22]$	0.988	-83.957	0.000	0.988	-83.957
4	$\alpha = [22, 21]$ $\beta = [0.1, 0.67]$	0.866	-83.724	0.000	0.886	-83.466
5	$\alpha = [0.8, 0.5]$ $\beta = [5, 22]$	0.860	-83.353	34.827	0.817	-84.009
6	$\alpha = [0.2, 0.4]$ $\beta = [1, 18]$	0.931	-70.692	35.754	0.821	-70.391